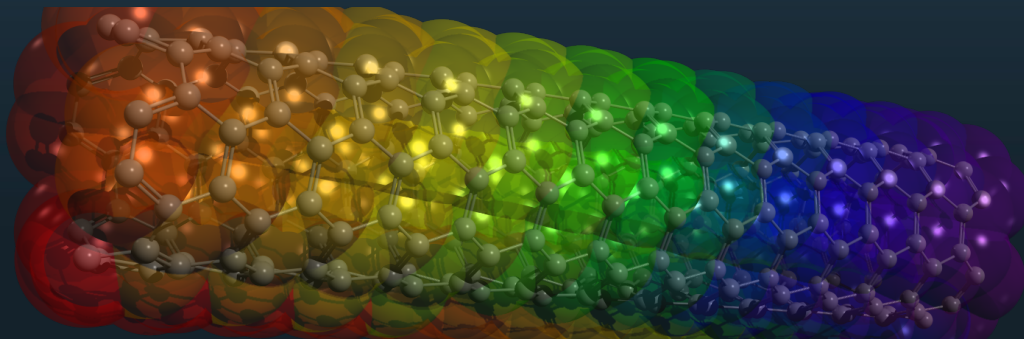# Chemical Databases and Open Chemistry on the Desktop

5th Meeting on US Government Chemical Databases & Open Chemistry
August 25, 2011

Dr. Marcus D. Hanwell
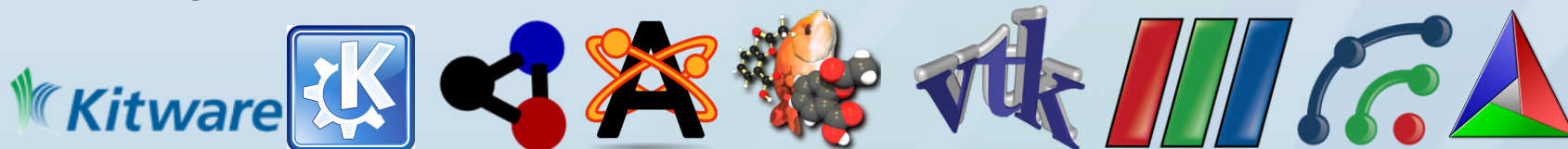marcus.hanwell@kitware.com

# Outline

- Background
- Opening up chemistry
- Workflows in computational chemistry
- Avogadro – chemical editor
- Databases on the desktop
- Quixote
- HPC resource integration
- Advanced visualization

**Kitware**

# My Background

- Ph.D. (Physics) – University of Sheffield
- Google Summer of Code – Avogadro
- Postdoc (Chemistry) – University of Pittsburgh
- R&D engineer – Kitware, Inc
- Passionate about physics, chemistry, and the growing need to improve computational tools
- See the need for powerful open source, cross platform frameworks and applications
- Develop(ed): Gentoo, KDE, Kalzium, Avogadro, Open Babel, VTK, ParaView, Titan, CMake
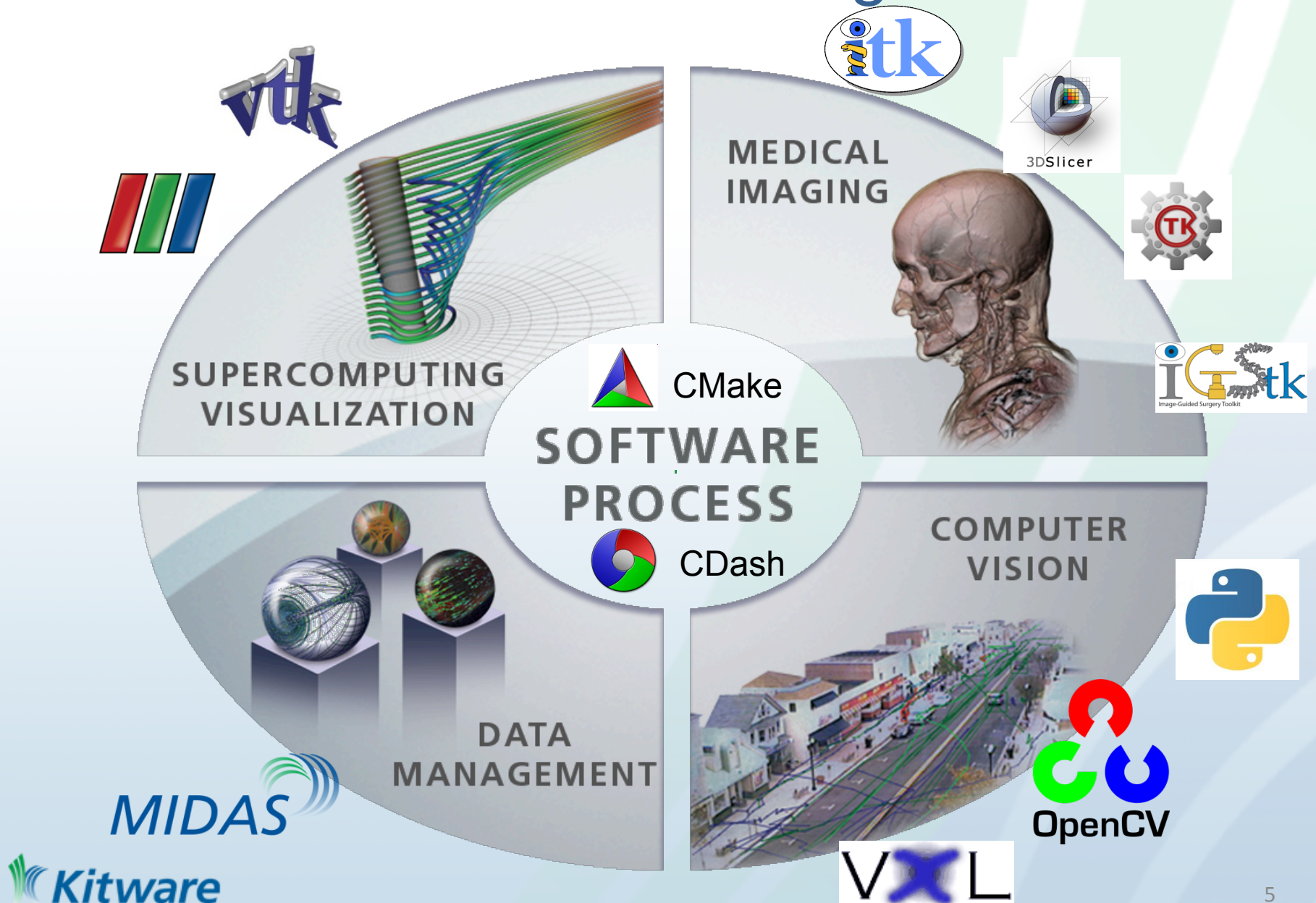
# Kitware

- Founded in 1998: 5 former GE Research employees

- 95 employees: 42% PhD

- Privately held, profitable from creation, no debt

- Rapidly Growing: >30% in 2010, 7M web-visitors/quarter

- Offices
  - Albany, NY
  - **Carrboro, NC**
  - Lyon, France
  - Bangalore, India



- 2011 Small Business Administration's Tibbetts Award

- HPCWire Readers and Editor's Choice

- Inc's 5000 List: 2008 to 2010

Kitware

# Kitware: Core Technologies
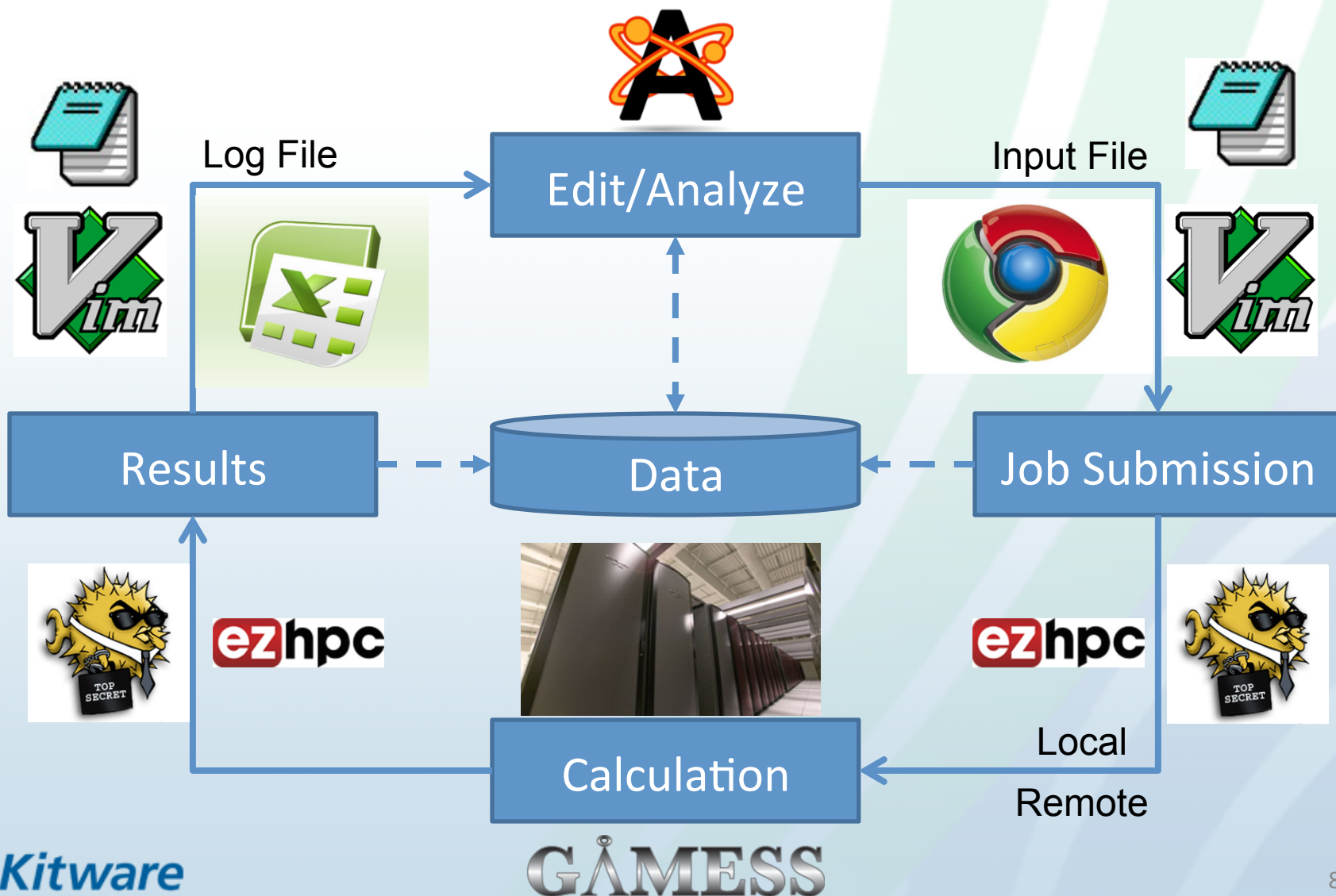
# Opening Up Chemistry

- Computational chemistry is currently one of the more closed sciences
- Lots of black box proprietary codes
  - Only a few have access to the code
  - Publishing results from black box codes
  - Many file formats in use, little agreement
- More papers should be including data
- Growing need for open standards

Kitware

# Movements for Open Chemistry

- Formed an "unorganization" – Blue Obelisk
  - Published first article in 2005
  - Open data, open standards and open source
  - Meet at ACS and other conferences when possible
  - Follow-up article currently in press

- Quixote collaboration more recently
  - Provide meaningful data storage and exchange
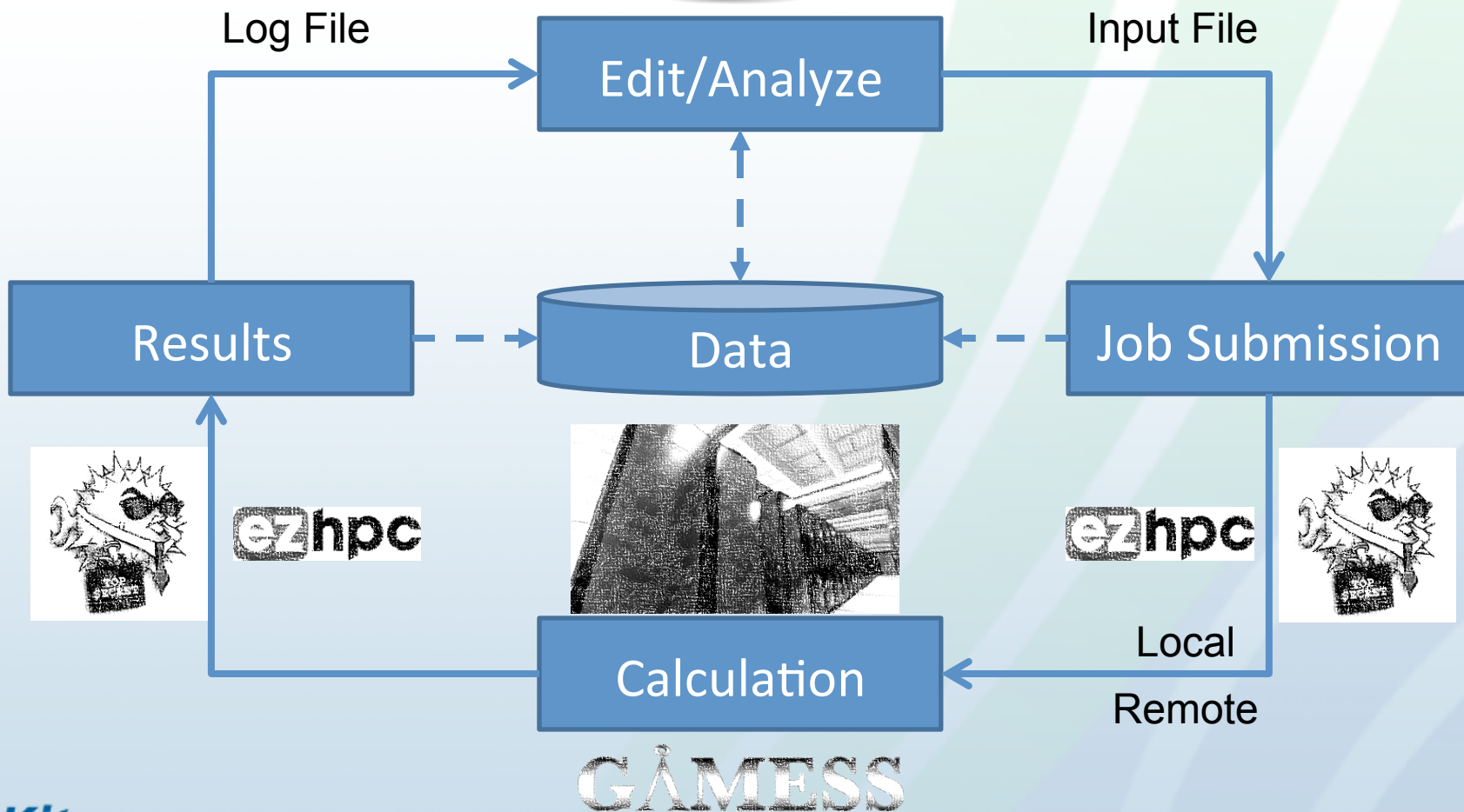  - Principally targeting computational chemistry

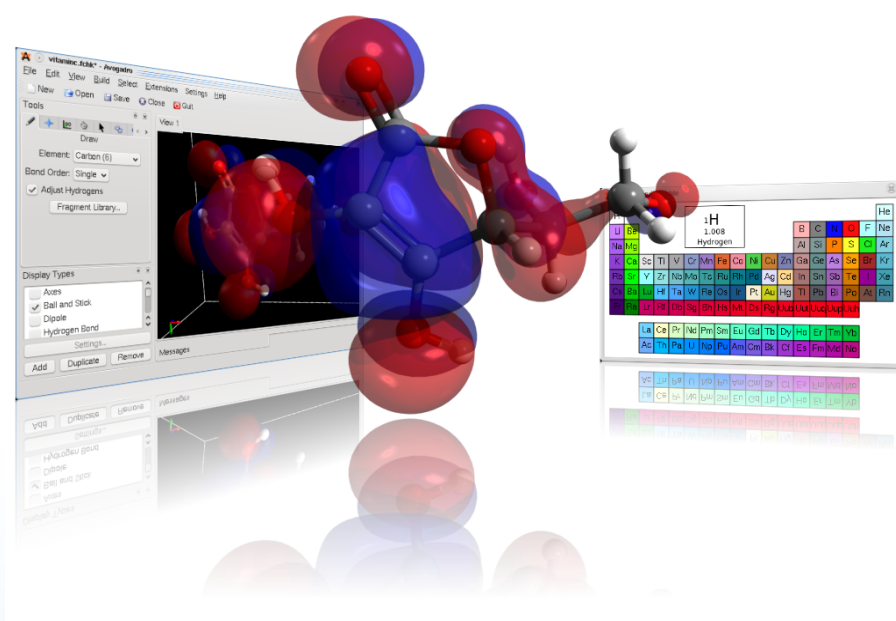**Kitware**

# Typical Chemistry Workflow

# Problem: Pretty Complex/Manual

- Most steps require user intervention
- Obtain starting structure (previous work, databases)
- Edit structure
- Write input file
- Move input file to cluster
- Submit to queue
- Wait for completion
- Retrieve input file
- Analyze output file
- Extract the relevant data, change formats
- Store results
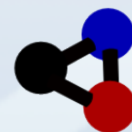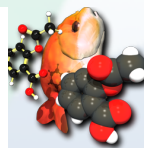- Repeat

**Kitware**

# Improved Chemistry Workflow

# Avogadro



- Project began 2006
- Split into library and application (plugin based)
- One of very few open source **editors**
- Designed to be extensible from the start
- Generate input & read output from many codes
- An active and growing community
- Chemistry needs a free, open framework

Kitware

# Avogadro's Roots

- Avogadro projected started in 2006
- First funded work in 2007 by Marcus Hanwell
  - Google Summer of Code student
  - Final year of Ph.D. spent the summer coding
  - Funded as part of KDE project – Kalzium editor
- Built on several other open source projects
  - Qt, Eigen, Open Babel, Blue Obelisk Data Repository
- Also uses open standards, e.g. OpenGL
- Cross platform, open source stack

# Avogadro Vital Statistics

- Supports Linux, Windows and Mac OS X
- Contributions from over 20 developers
- Over 180,000 downloads over 4 years
- Translated into 19 languages
- Used by Kalzium for molecular editor
- Featured by Trolltech/Nokia,
  - Qt in use
  - Qt ambassador program

# Desktop Database

- Use of "document store" NoSQL
  - Doesn't force too much structure
    - Some entries have experimental data available
    - Some have computational jobs
  - Employ a "pile of stuff" approach
    - Can store both source and derived data
    - Calculate identifiers, QSAR properties, etc
- MongoDB is a scalable, open solution
  - Proven scaling with large web applications

*Kitware*

# Chemistry Data Explorer

- Qt application
- Connects to local or remote database
- Uses VTK for visual data exploration
- Can ingest new data
  - Uses Open Babel to generate descriptors
  - Standard InChi, SMILES, molecular weight
  - More could be added
    - All derived from files stored in the database

**Kitware**

# Chemistry Data Explorer

# Database Interaction on the Web

- Avogadro directly accesses some (read-only) public databases:
  - PDB, NIH "fetch by name"
  - Resolve structure to common name using CIR
  - More could be added
- ChemData also uses NIH CIR for data
- Quixote aims to support both public and private sharing models – open framework

**Kitware**

# Quixote Architecture

# OpenQube – Quantum Data

- Reads in key quantum data
  - Basis set used in calculation
  - Eigenvectors for molecular orbitals
  - Density matrix for electron density
  - Standard geometry
- Multithreaded calculation
  - Produce regular grids of scalar data
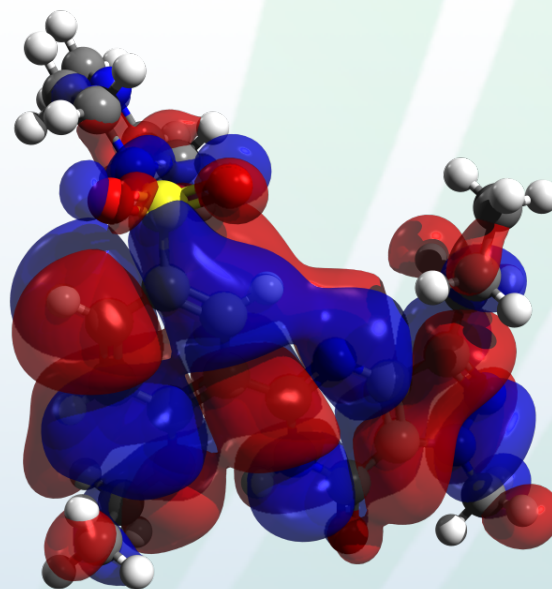  - Molecular orbitals, electron density…

*Kitware*

# Molecular Orbitals and Electron Density

- Quantum files store basis sets and matrices

$$GTO = ce^{-\alpha r^2}$$

$$\phi_i = \sum_\mu c_{\mu i} \phi_\mu$$

$$\rho(r) = \sum_\mu \sum_\nu P_{\mu\nu} \phi_\mu \phi_\nu$$

- Using these equations, and the supplied matrices – calculate cubes

**Kitware**

# Calling Stand Alone Programs

- Many already supported:
  - GAMESS, GAMESS-UK, Molpro, Q-Chem, MOPAC, NWChem, Gaussian, Dalton
  - Easy to add more
- Some codes writing Avogadro based custom applications,
  - Q-Chem, Molpro…
- DLPOLY author approached me:
  - Open sourced DLPOLY2, want a GUI

*Kitware*

# Job Submission & Management

- Take input file, submit to queue, monitor, retrieve, repeat
- System tray resident Qt application
  - Manage both local and remote jobs
- Interest from developers
  - Use in other applications
  - Share development/maintenance burden

**Kitware**

# Open in Avogadro When Complete

# Advanced Visualization: VTK

- New Avogadro plugin:
  - Takes volumetric data from Avogadro
  - Uses GPU accelerated rendering in VTK
- Excitement from many in the community
- Several groups interested in collaborating
- Google Summer of Code project
- Leverage significant capabilities in VTK

# Volume Rendered With Contours

# Electron Density Volume Render

# Electron Density Ray Tracing

# Conclusions

- There is still a lot of work to do

- Open databases are of critical importance

- Need tools to make retrieving and depositing data easier

- Improved data exchange is essential to improve reproducibility in chemistry

- Create shared collaboration platforms
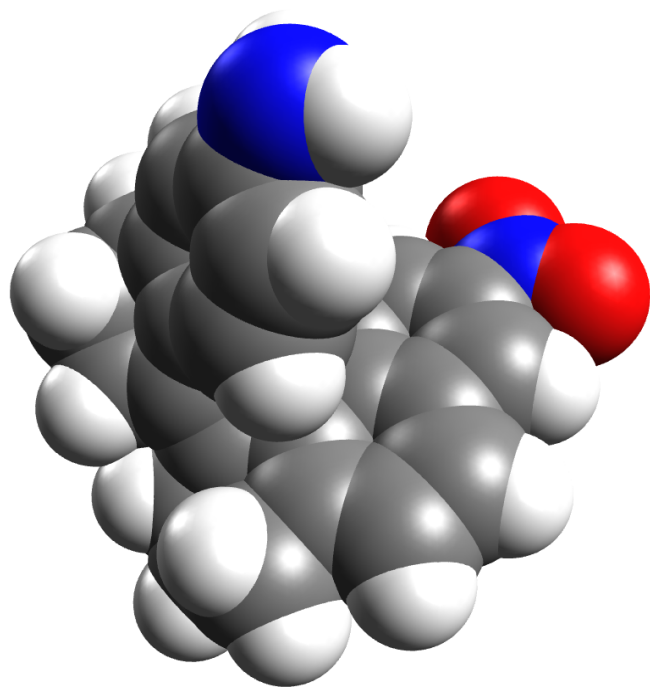  - Deliver improved workflows, enable research

**Kitware**

# Extra Background Slides

- Additional visualization and background slides
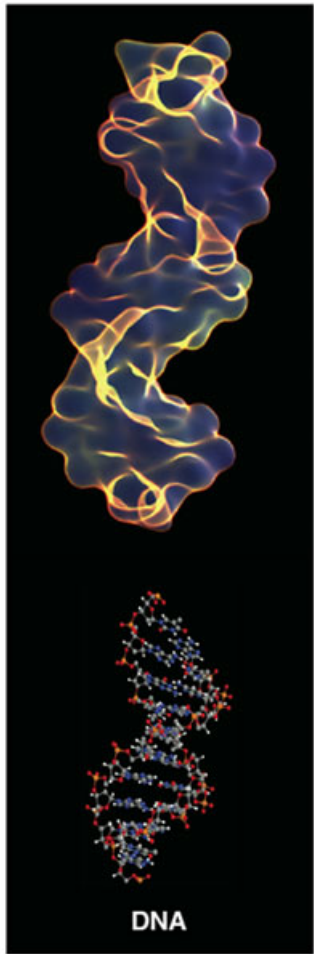
# Standard Representations
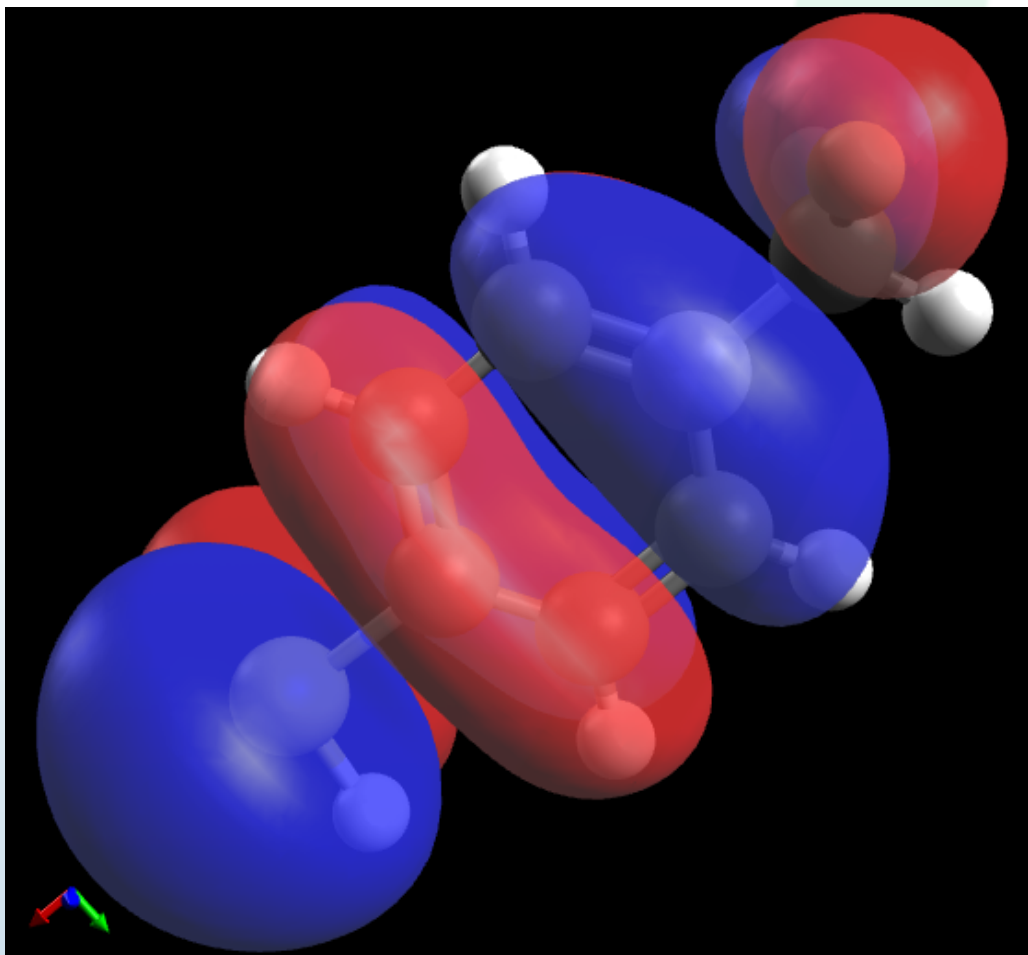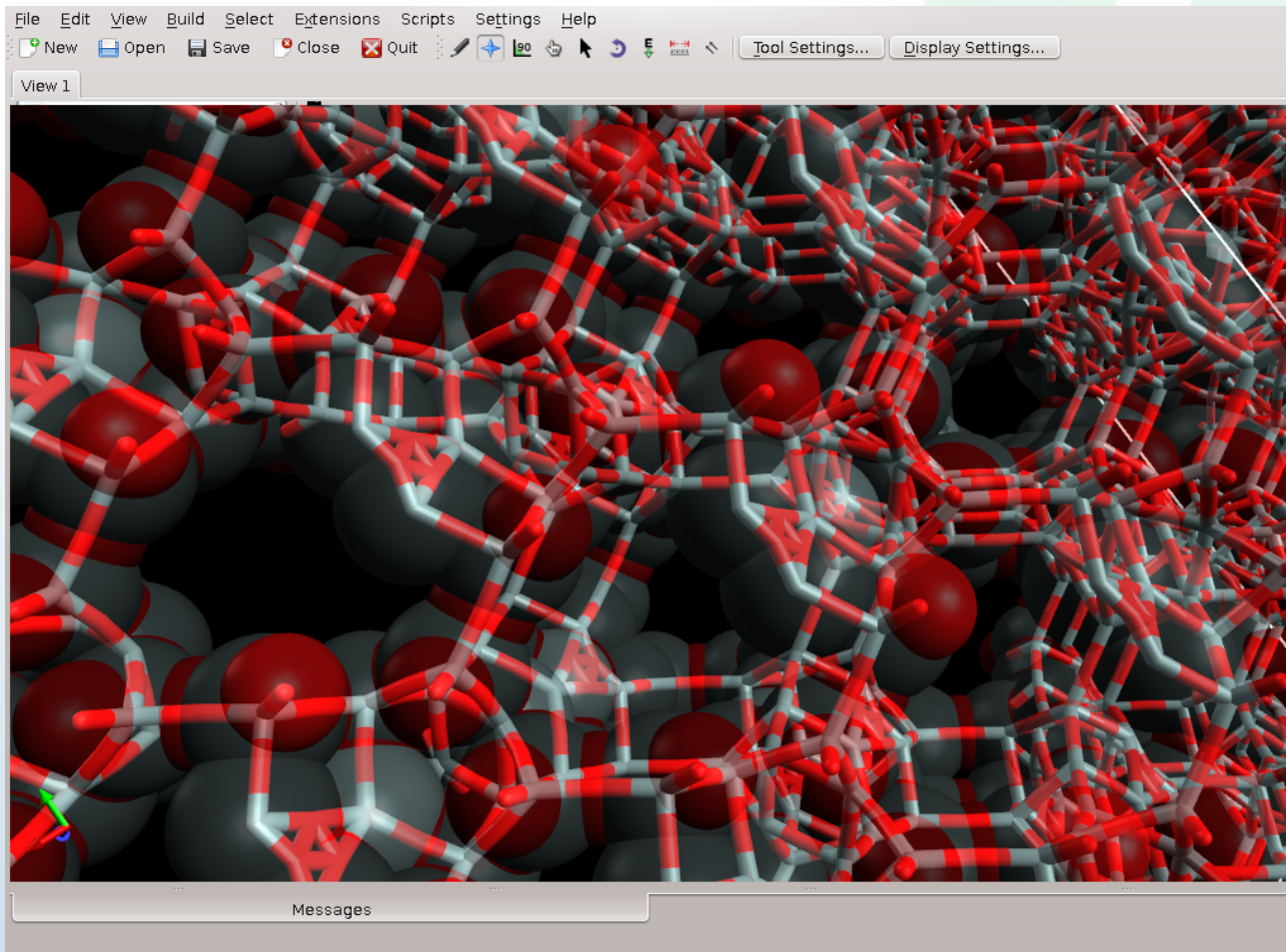
# Standard Representations
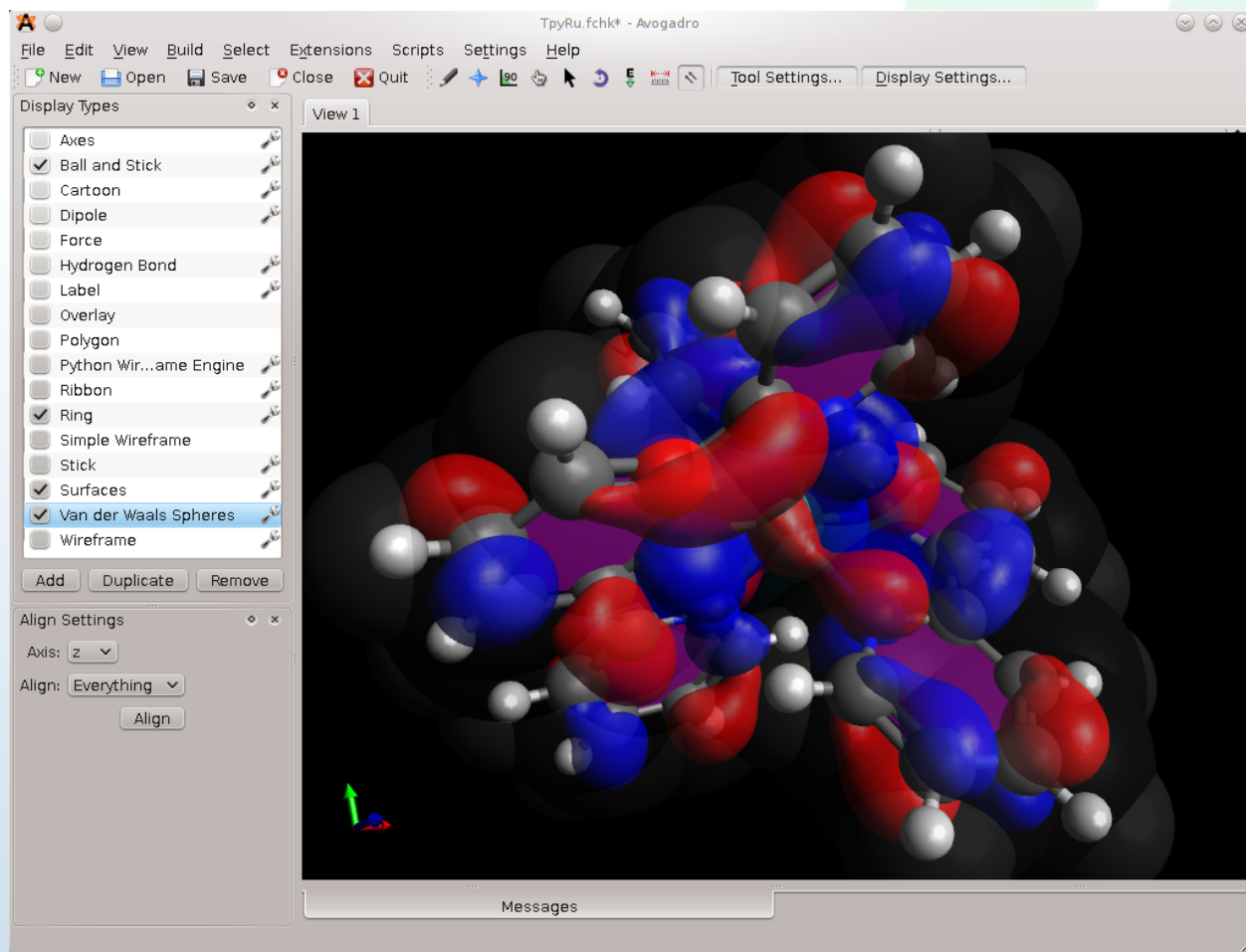
# Biomolecules

# Nanomaterials

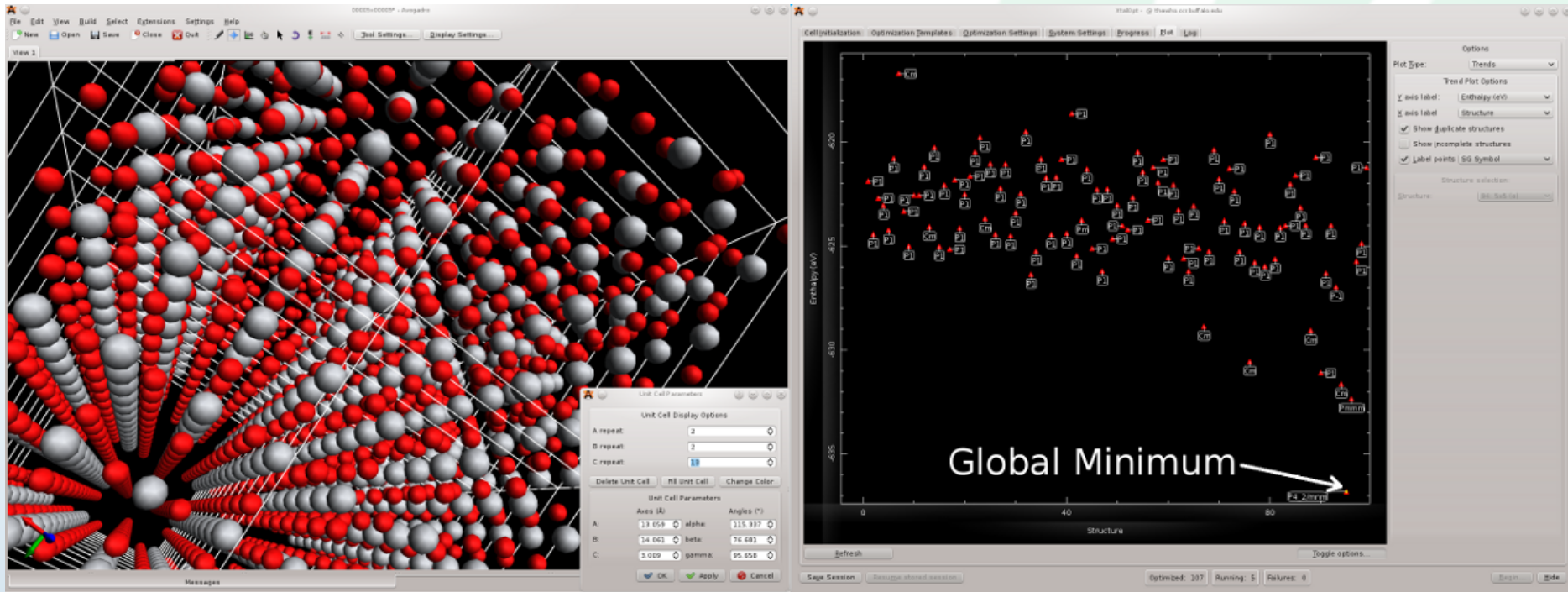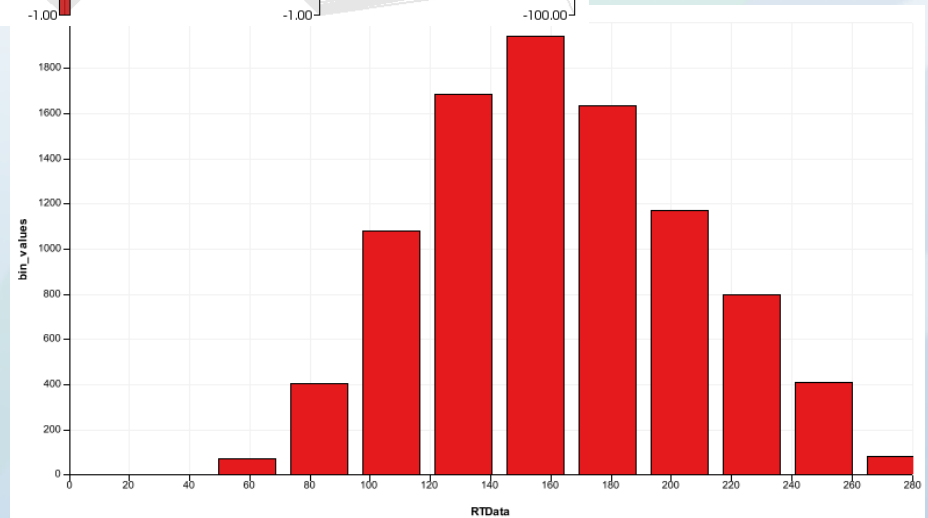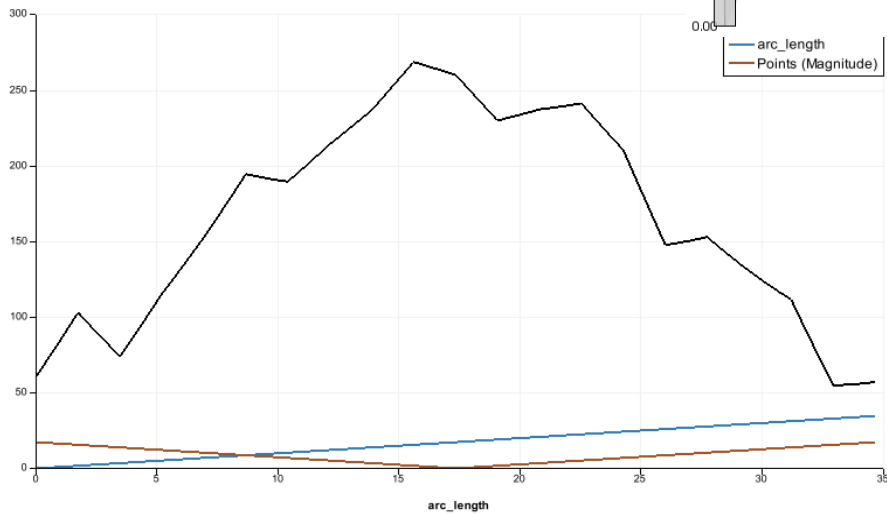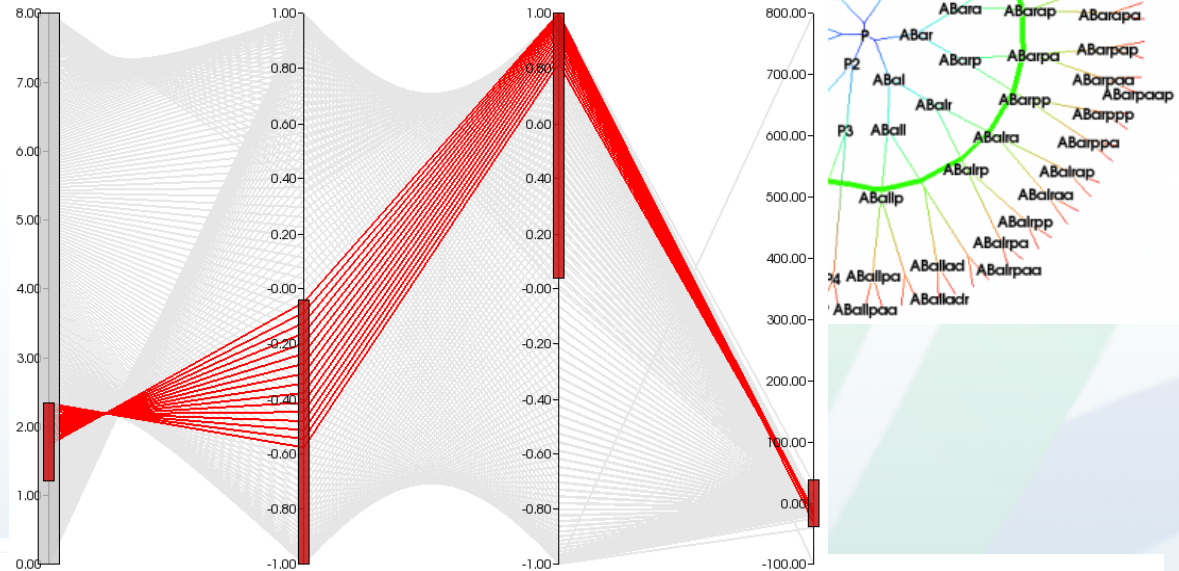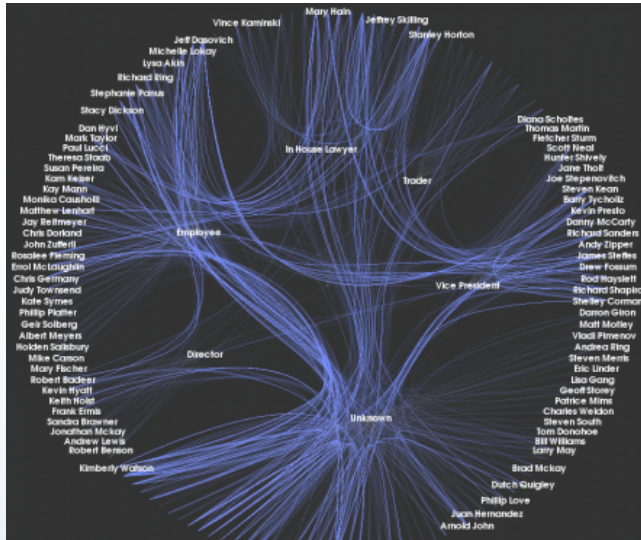# Simplified Views

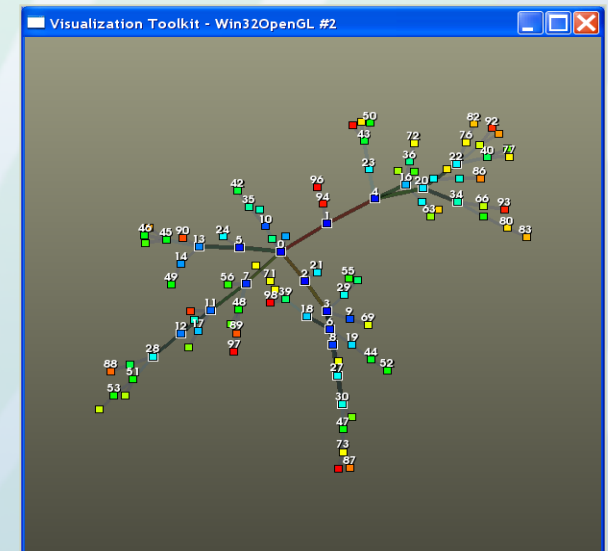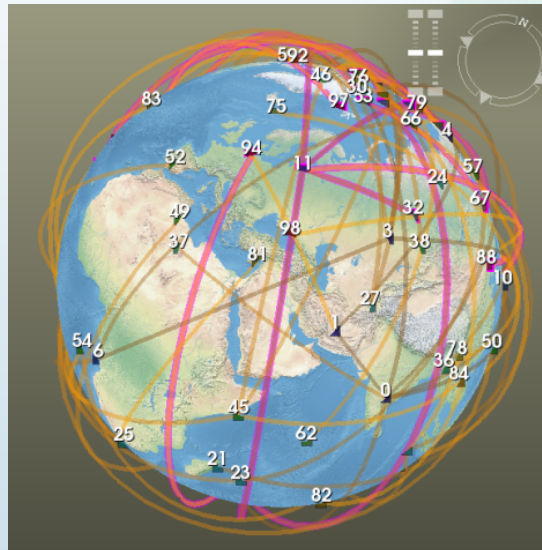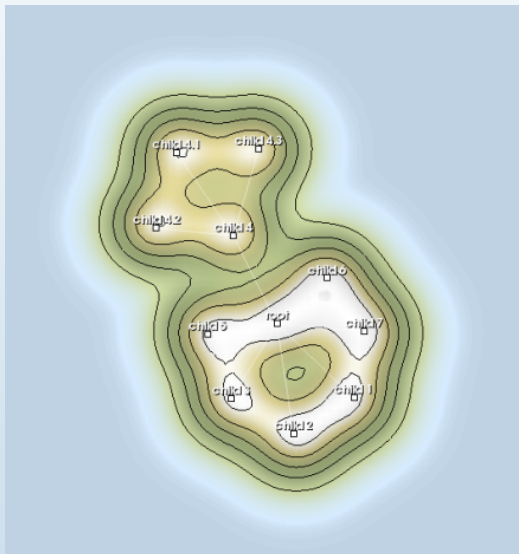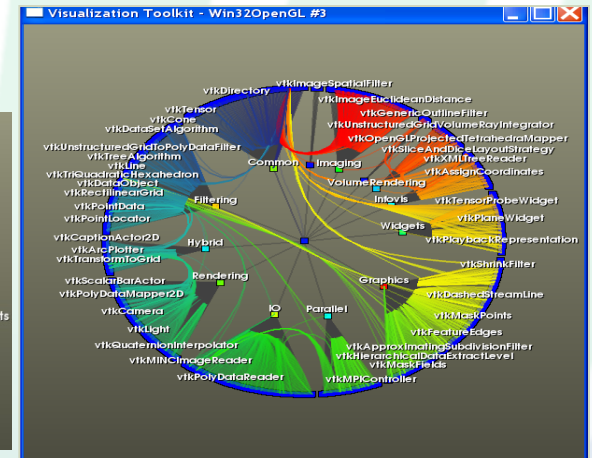
DNA

# Volumetric Data: Molecular Orbitals

# Periodic Systems

# Hybrid Views: CPK + MO + Ball & Stick

# Linked Views of Live Data



Global Minimum

# 2D: Graphs and Charts

# Informatics

# 3D Interaction Widgets

# VTK: The Toolkit

- Collection of C++ libraries
  - Leveraged by many applications
  - Divided into logical areas, e.g.
    - Filtering – data processing in visualization pipeline
    - InfoVis – informatics visualization
    - Widgets – 3D interaction widgets
    - VolumeRendering – 3D volume rendering
- Cross platform, using OpenGL
- Wrapped in Python, Tcl and Java

# VTK Development Team

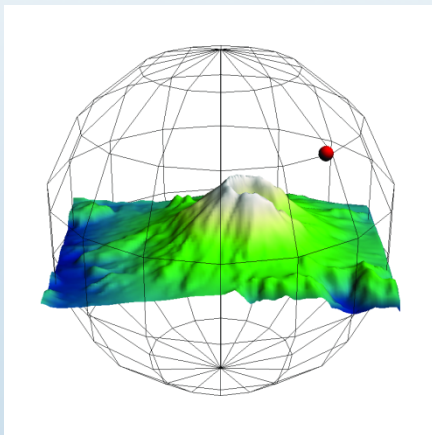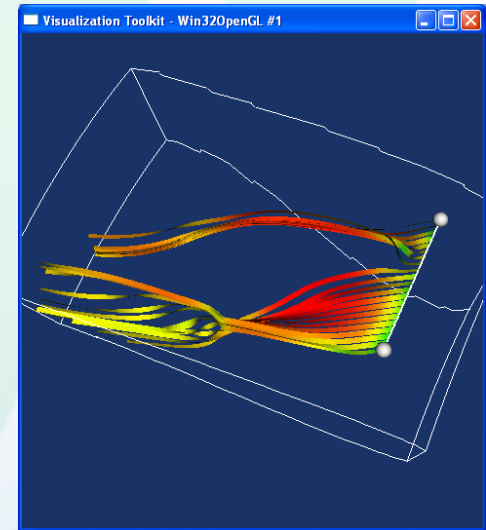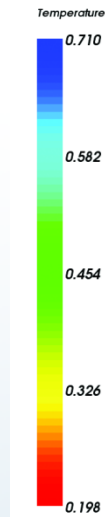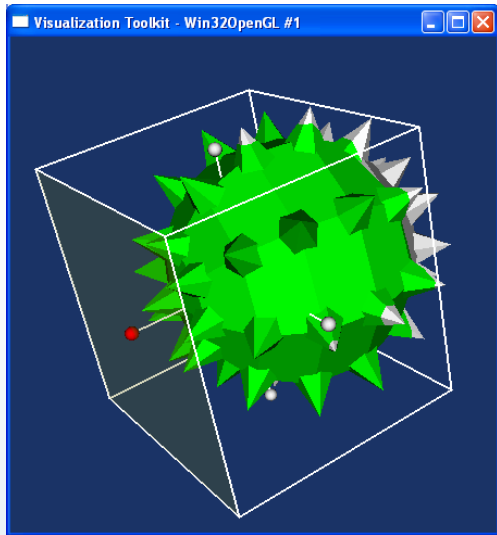- From Ohloh: **Very large, active development team:** Over the past twelve months, **100 developers** contributed new code to VTK. This is one of the largest open-source teams in the world, and is in the **top 2%** of all project teams on Ohloh.

PITTSBURGH SUPERCOMPUTING CENTER

eDF

CSCS
Swiss National Supercomputing Centre

Los Alamos
NATIONAL LABORATORY
EST. 1943

CD-adapco

ARL
U.S. Army Research Laboratory

GE

TUDelft
Delft University of Technology

Kitware

THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

Sandia National Laboratories

INDIANA UNIVERSITY

SCI
www.sci.utah.edu

GEORGETOWN UNIVERSITY

cea

JOHNS HOPKINS UNIVERSITY

MIT
Massachusetts Institute of Technology

Rensselaer

and many others...

Kitware

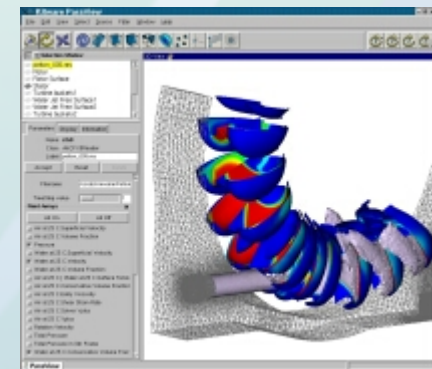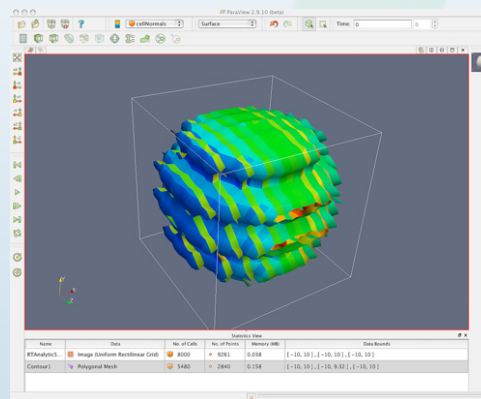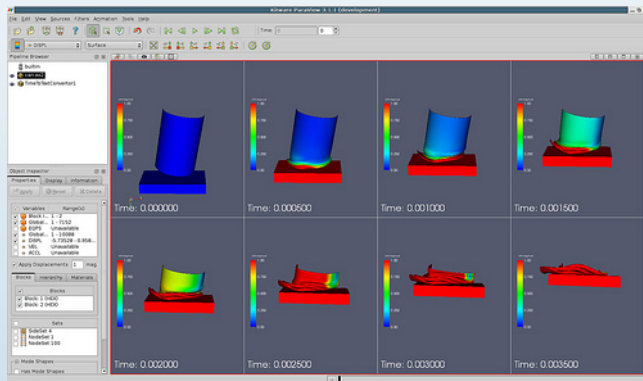# ParaView



- Parallel visualization application
- Open source, BSD licensed
- Turn-key application wrapper around VTK
- Parallel data processing **and** rendering
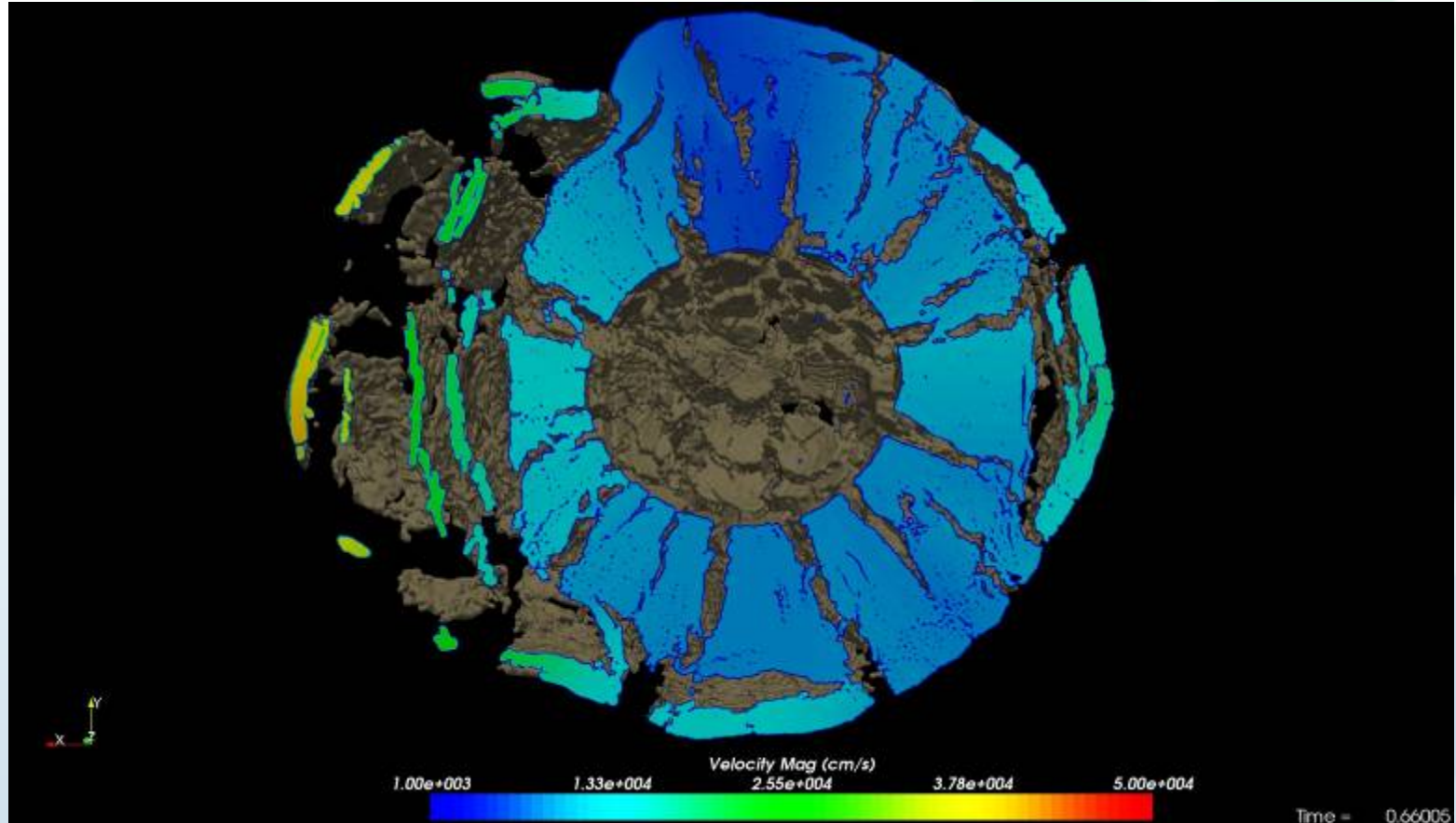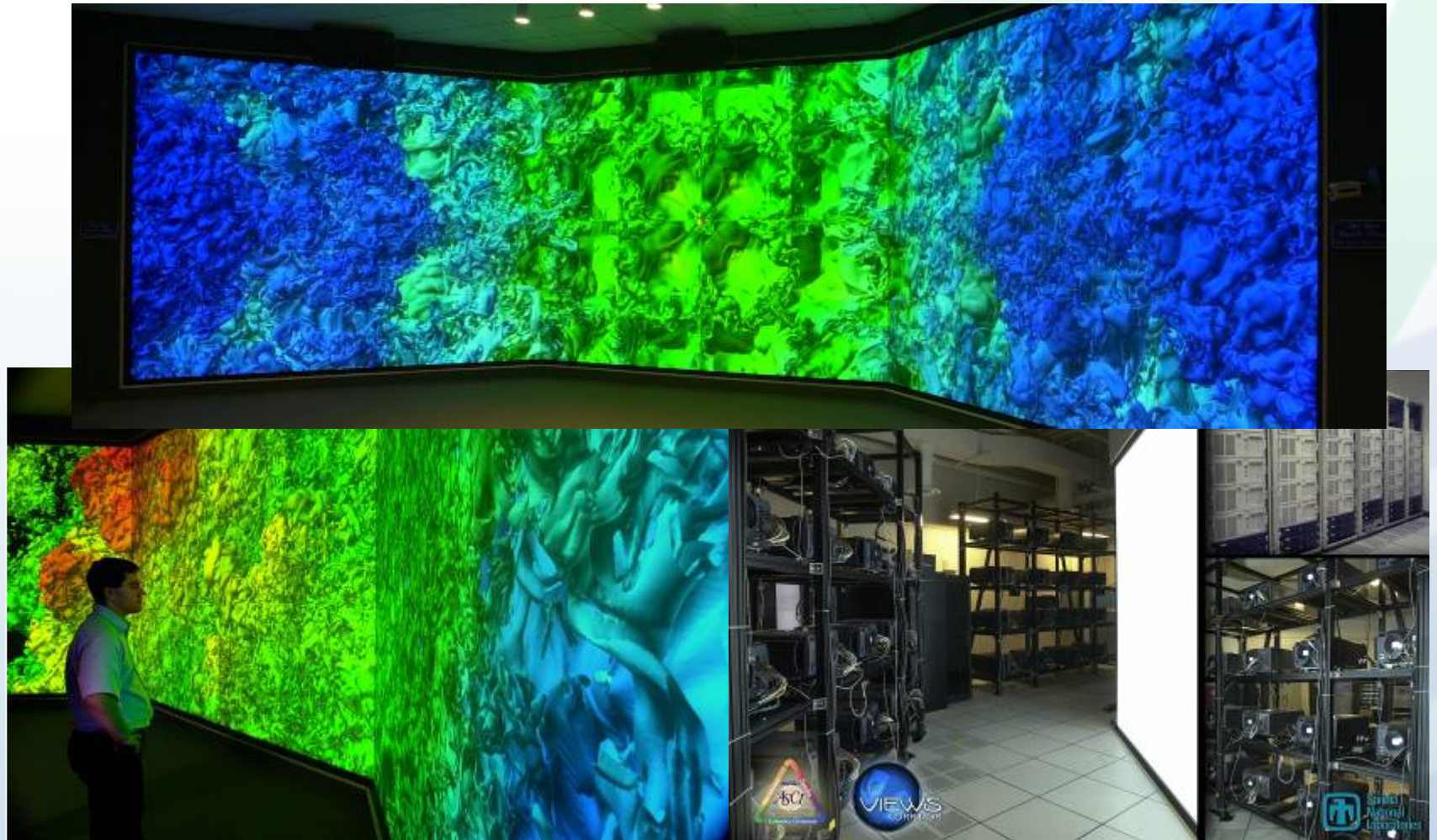
# Large Data Visualization

- BlueGene/L at LLNL
  - 65,536 compute nodes (32 bit PPC)
  - 1,024 I/O nodes (32 bit PPC)
  - 512 MB of RAM per node
- Sandia Red Storm
  - 12,960 compute nodes (AMD Opteron dual)
  - 640 service and I/O nodes
  - 40 TB of DDR RAM per node





Kitware
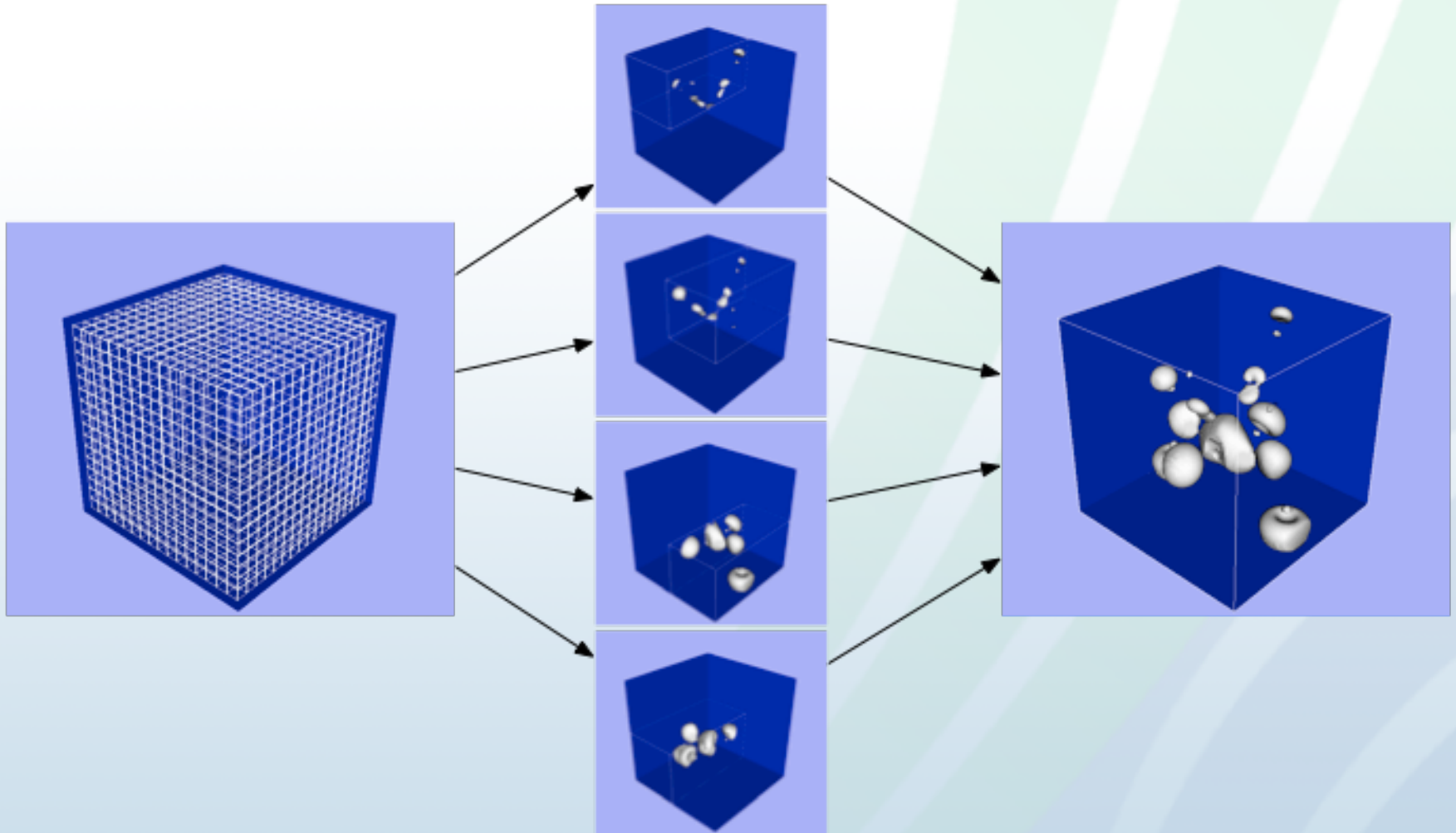
# 1 Billion Cell Asteroid Simulation
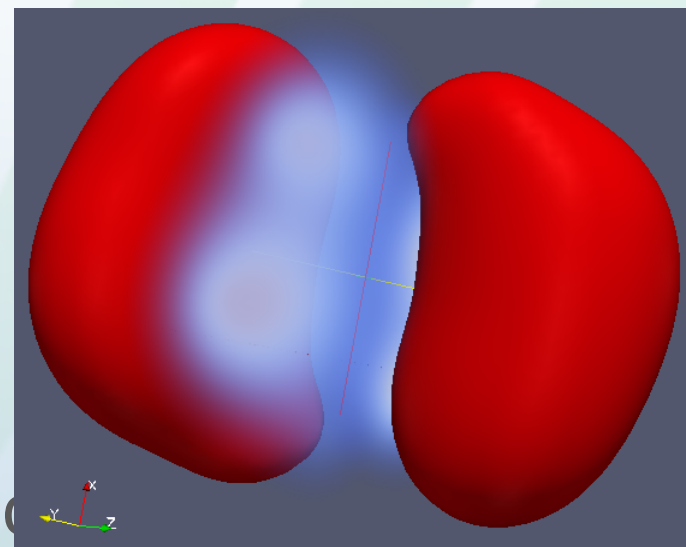


Kitware

# Tiled Displays

# Parallel Processing/Rendering
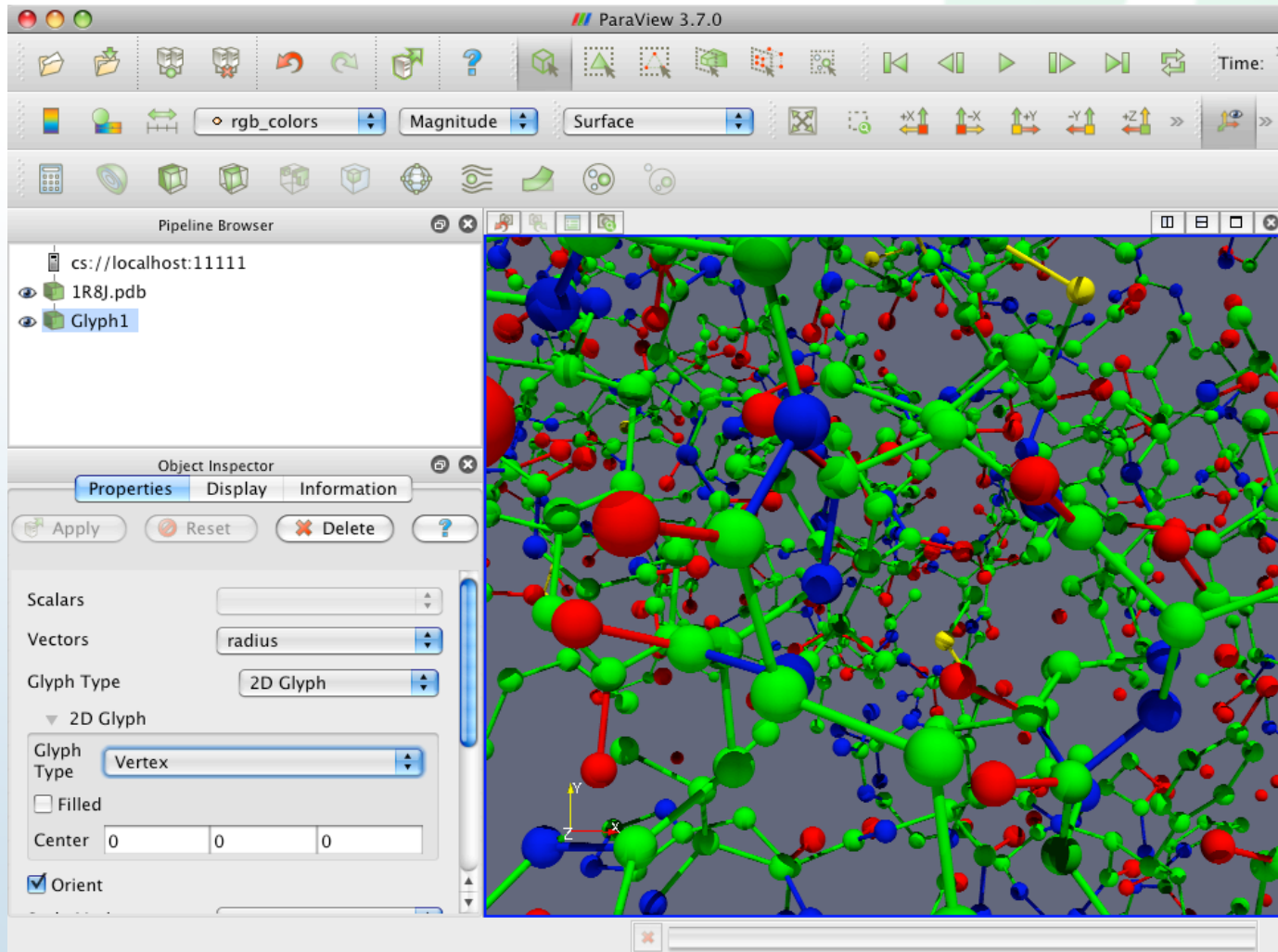
# 3D Chemistry Visualization

- Some existing features specific to chemistry
  - Gaussian cube, PDB, and a few others
- Excellent handling of volumetric data:
  - Marching cubes
  - Volume rendering
  - Contouring
- Advanced rendering:
  - Point sprites
  - Manta – real time ray trac

Kitware

# Titan: VTK and Informatics

- Led by Sandia National Laboratories
- Substantial expansion of VTK:
  - Informatics & analysis
- Actively developed, growing feature set
- Improved 2D rendering and API
- Database connectivity, client-server, pipeline based approach
- Uses web technologies such as ProtoViz
- Scalable, interactive infoviz

# Manta: Real Time Ray Tracing

# New Frontiers

- New work porting VTK
  - Use C++ as the common core
    - iOS port in the early stages
    - Android port
  - Use OpenGL ES 2.0 – new rendering code
- Also ParaViewWeb – delivering over web
  - Use image delivery and rendering on server
  - Also using WebGL for rendering (optionally)

# Future Directions

- VTK modularization (in progress)
  - Developing more agile build systems
  - Automating more with CMake
- Using Git more fully to improve stability
  - Use of master and next
  - Topic branches - merge when ready
- Code review using Gerrit
  - Integration with continuous integration
  - Test before merge

*Kitware*